

# Case-Based Multiagent Reinforcement Learning: Cases as Heuristics for Selection of Actions

Reinaldo A. C. Bianchi<sup>1,2</sup> and Ramón López de Mántaras<sup>2</sup>

**Abstract.** This work presents a new approach that allows the use of cases in a case base as heuristics to speed up Multiagent Reinforcement Learning algorithms, combining Case-Based Reasoning (CBR) and Multiagent Reinforcement Learning (MRL) techniques. This approach, called Case-Based Heuristically Accelerated Multiagent Reinforcement Learning (CB-HAMRL), builds upon an emerging technique, Heuristic Accelerated Reinforcement Learning (HARL), in which RL methods are accelerated by making use of heuristic information. CB-HAMRL is a subset of MRL that makes use of a heuristic function  $\mathcal{H}$  derived from a case base, in a Case-Based Reasoning manner. An algorithm that incorporates CBR techniques into the Heuristically Accelerated Minimax-Q is also proposed and a set of empirical evaluations were conducted in a simulator for the Littman's robot soccer domain, comparing the three solutions for this problem: MRL, HAMRL and CB-HAMRL. Experimental results show that using CB-HAMRL, the agents learn faster than using RL or HAMRL methods.

## 1 Introduction

Heuristic Accelerated Reinforcement Learning (HARL) [6] is an emerging technique in which Reinforcement Learning (RL) methods are sped up by making use of a conveniently chosen heuristic function, which is used for selecting appropriate actions to perform in order to guide exploration during the learning process. HARL techniques are very attractive: as RL, they are based on firm theoretical foundations. As the heuristic function is used only in the choice of the action to be taken, many of the conclusions obtained for RL remain valid for HARL algorithms, such as the guarantee of convergence to equilibrium in the limit – given that some predefined conditions are satisfied – and the definition of an upper bound for the error [6].

Although several methods have been successfully applied for defining the heuristic function, a very interesting option has only recently been explored: the reuse of previously learned policies, using a Case-Based Reasoning approach [8]. This paper investigates the combination of Case-Based Reasoning (CBR) and Multiagent Heuristically Accelerated Reinforcement Learning (HAMRL) [7] techniques, with the goal of speeding up MRL algorithms by using previous domain knowledge, stored as a case base. To do so, we propose a new algorithm, the Case-Based Heuristically Accelerated Minimax-Q (CB-HAMMQ), which incorporates Case-Based Reasoning techniques into an existing HAMRL algorithm, the Heuristically Accelerated Minimax-Q (HAMMQ).

Soccer competitions, such as RoboCup, have been proven to be an important challenge domain for research, and one where RL tech-

niques have been widely used. The application domain of this paper is a simulator for the robot soccer domain that extends the one proposed by Littman [20], called “Expanded Littman's Soccer”. Nevertheless, the technique proposed in this work is domain independent.

The paper is organized as follows: section 2 briefly reviews the Multiagent Reinforcement Learning problem, describes the HAMRL approach and the HAMMQ algorithm, while section 3 describes Case-Based Reasoning. Section 4 shows how to incorporate CBR techniques into HAMRL algorithms, in a modified formulation of the HAMMQ algorithm. Section 5 describes the Robotic Soccer domain used in the experiments, presents the experiments performed, and shows the results obtained. Finally, Section 6 provides our conclusions.

## 2 Heuristic Accelerated Multiagent Reinforcement Learning

Systems where multiple agents compete among themselves to accomplish their tasks can be modeled as a discrete time, finite state, finite action Markov Game (MG) – also known as Stochastic Game (SG). The goal of an agent in a MRL problem is to learn an optimal policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}_1 \times \dots \times \mathcal{A}_k$  that maps the current state  $s$  into a desirable action(s)  $a$  to be performed in  $s$ , from any starting state. In MRL, this policy is learned through trial-and-error interactions of the agent with its environment: on each interaction step the agent senses the current state  $s$  of the environment, chooses an action  $a$  to perform, executes this action, altering the state  $s$  of the environment, and receives a scalar reinforcement signal  $r$  (a reward or penalty).

This paper considers a well-studied specialization of MGs in which there are only two players, called agent and opponent, having opposite goals. Such specialization, called a zero-sum Markov Game (ZSMG) [20], allows the definition of only one reward function that the learning agent tries to maximize while the opponent tries to minimize. A two player ZSMG is defined by the quintuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{O}, T, \mathcal{R} \rangle$ , where:

- $\mathcal{S}$ : a finite set of environment states.
- $\mathcal{A}$ : a finite set of actions that the agent can perform.
- $\mathcal{O}$ : a finite set of actions that the opponent can perform.
- $T : \mathcal{S} \times \mathcal{A} \times \mathcal{O} \rightarrow \Pi(\mathcal{S})$ : the state transition function, where  $\Pi(\mathcal{S})$  is a probability distribution over the set of states  $\mathcal{S}$ .  $T(s, a, o, s')$  defines a probability of transition from state  $s$  to state  $s'$  (at a time  $t + 1$ ) when the learning agent executes action  $a$  and the opponent performs action  $o$ .
- $\mathcal{R} : \mathcal{S} \times \mathcal{A} \times \mathcal{O} \rightarrow \mathbb{R}$ : the reward function that specifies the reward received by the agent when it executes action  $a$  and its opponent performs action  $o$ , in state  $s$ .

<sup>1</sup> Centro Universitário da FEI, São Bernardo do Campo, Brazil.

<sup>2</sup> Artificial Intelligence Research Institute (IIIA-CSIC), Bellaterra, Spain.

To solve a ZSMG, Littman [20] proposed the use of a strategy similar to Minimax for choosing an action in the Q-Learning algorithm, the Minimax-Q algorithm, which works in the same way as Q-Learning does. The action-value function of an action  $a$  in a state  $s$  when the opponent takes an action  $o$  is can be computed iteratively by:

$$\hat{Q}_{t+1}(s, a, o) \leftarrow \hat{Q}_t(s, a, o) + \alpha [r(s, a, o) + \gamma V_t(s') - \hat{Q}_t(s, a, o)], \quad (1)$$

where  $\alpha$  is the learning rate,  $\gamma$  is the discount factor and the value  $V_t(s)$  of a state can be computed using the following equation:

$$V(s) = \max_{\pi \in \Pi(\mathcal{A})} \min_{o \in \mathcal{O}} \sum_{a \in \mathcal{A}} Q(s, a, o) \pi_a, \quad (2)$$

where the agent's policy  $\pi$  is a probability distribution over actions, and  $\pi_a$  is the probability of taking the action  $a$  against the opponent's action  $o$ . In an Alternating Markov Game (AMG), where two players take their actions in consecutive turns, the policy becomes deterministic and Equation 2 can be simplified:

$$V(s) = \max_{a \in \mathcal{A}} \min_{o \in \mathcal{O}} Q(s, a, o). \quad (3)$$

Formally, a Heuristically Accelerated Multiagent Reinforcement Learning (HAMRL) algorithm is a way to solve a MG problem with explicit use of a heuristic function  $\mathcal{H} : \mathcal{S} \times \mathcal{A} \times \mathcal{O} \rightarrow \mathbb{R}$  to influence the choice of actions during the learning process.  $H(s, a, o)$  defines a heuristic that indicates the desirability of performing action  $a$  when the agent is in state  $s$  and the opponent executes action  $o$ .

The first HAMRL algorithm proposed was the Heuristically Accelerated Minimax Q (HAMMQ) [7], as an extension of the Minimax-Q algorithm. The only difference between them is that in the HAMMQ the heuristic function is used in the action choice rule, which defines which action  $a_t$  must be executed when the agent is in state  $s_t$ . The action choice rule used in the HAMMQ is a modification of the standard  $\epsilon - Greedy$  rule used in Minimax-Q, to include the heuristic function:

$$\pi(s) = \begin{cases} \arg \max_a \min_o [\hat{Q}(s, a, o) + \xi H_t(s, a, o)] & \text{if } q \leq p, \\ a_{random} & \text{otherwise,} \end{cases} \quad (4)$$

where  $\mathcal{H} : \mathcal{S} \times \mathcal{A} \times \mathcal{O} \rightarrow \mathbb{R}$  is the heuristic function,  $q$  is a random value uniformly distributed over  $[0, 1]$  and  $0 \leq p \leq 1$  is a parameter that defines the exploration/exploitation tradeoff. The subscript  $t$  indicates that the heuristic function can be non-stationary (it can be computed only once, or be continually recomputed) and  $0 \leq \xi \leq 1$  is a real variable used to weight the influence of the heuristic.

As a general rule, the value of  $H_t(s, a, o)$  used in HAMMQ should be higher than the variation among the  $\hat{Q}(s, a, o)$  values for the same  $s \in \mathcal{S}$ ,  $o \in \mathcal{O}$ , in such a way that it can influence the choice of actions, and it should be as low as possible in order to minimize the error. It can be defined as:

$$H(s, a, o) = \begin{cases} \max_i \hat{Q}(s, i, o) - \hat{Q}(s, a, o) + \eta & \text{if } a = \pi^H(s), \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$

where  $\eta$  is a small real value (usually 1) and  $\pi^H(s)$  is the action suggested by the heuristic policy. Convergence of this algorithm was presented by Bianchi, Ribeiro and Costa [7], together with the definition of an upper bound for the error. The complete HAMMQ algorithm is presented in Table 1.

**Table 1.** The HAMMQ algorithm.

---

|   |
|---|
| Initialize $\hat{Q}_t(s, a, o)$ and $H_t(s, a, o)$ arbitrarily. |
| Repeat (for each episode):                                      |
| Initialize $s$ .  |
| Repeat (for each step):   |
| Update the values of $H_t(s, a, o)$ as desired.                 |
| Select an action $a$ using Equation 4.                          |
| Execute the action $a$ , observe $r(s, a, o)$ , $s'$ .          |
| Update the values of $Q(s, a, o)$ according to Equation 1.      |
| $s \leftarrow s'$ .   |
| Until $s$ is terminal.  |
| Until some stopping criterion is reached.                       |

---

Despite the fact that RL is a method that has been traditionally applied in the Robotic Soccer domain, only recently have HARL methods been used in this domain. Bianchi, Ribeiro and Costa [7] investigated the use of a HAMRL algorithm in a simplified simulator for the robot soccer domain and Celiberto *et al.* [10] studied the use of the HARL algorithms to speed up learning in the RoboCup 2D Simulation domain. The heuristic used in both of these papers were very simple ones: in the first paper the heuristic was ‘if the agent is with the ball, go to the opponent’s goal’, and in the second paper it was simply ‘go to the ball’.

### 3 Case-Based Reasoning

Humans frequently try to solve a new problem by remembering a previous similar situation, reasoning about it, and then reusing knowledge of that situation to solve the new problem. Case-based reasoning (CBR) [1, 22] uses knowledge of previous situations (cases) to solve new problems, by finding a similar past case and reusing it in the new problem situation. In the CBR approach, a case usually describes a problem and its solution, i.e., the state of the world in a given instant and the sequence of actions to perform to solve that problem.

According to López de Mántaras *et al.* [22], solving a problem by CBR involves “obtaining a problem description, measuring the similarity of the current problem to previous problems stored in a case base with their known solutions, retrieving one or more similar cases, and attempting to reuse the solution of the retrieved case(s), possibly after adapting it to account for differences in problem descriptions”. Other steps that are usually found in CBR systems are the evaluation of the proposed solution, the revision of the solution, if required in light of its evaluation, and the retention (learning) of a new case, if the system has learned to solve a new problem.

The case definition used in this work is the one proposed in Ros [25] and Ros *et al.* [26], which is composed of three parts: the problem description ( $P$ ), the solution description ( $A$ ) and the case scope ( $K$ ), and it is formally described as a 3-tuple:

$$case = (P, A, K). \quad (6)$$

The problem description  $P$  corresponds to the situation in which the case can be used. For example, for a Robotic Soccer problem, the description of a case can include the robot position, the ball’s position and the positions of the other robots in the game. For a game with  $n$  robots (teammates and opponents),  $P$  can be:

$$P = \{x_B, y_B, x_{R_1}, y_{R_1}, \dots, x_{R_n}, y_{R_n}\}. \quad (7)$$

The solution description is composed by the sequence of actions that each robot must perform to solve the problem, and can be defined

as:

$$A = \{R_1 : [a_{11}, a_{12}, \dots, a_{1p_1}], \dots, R_m : [a_{m1}, a_{m2}, \dots, a_{mp_m}]\},$$

where  $m$  is the number of robots in the team,  $a_{ij}$  is an individual or joint action that robot  $R_i$  must perform and  $p_i$  corresponds the number of actions the robot  $R_i$  performs.

The case scope defines the applicability boundaries of the cases, to be used in the retrieval step. For example, Ros [25] defines it as “the regions of the field within which the ball and the opponents should be positioned in order to retrieve that case”. In the case of a robot soccer problem,  $K$  can be represented as circles or ellipsoids centered on the ball’s and opponents’ positions indicated in the problem description. It can be defined as:

$$K = \{\tau_B, \tau_{R_1}, \dots, \tau_{R_n}\}, \quad (8)$$

where  $\tau_B$  is the radius of the region around the ball and  $\tau_{R_1} \dots \tau_{R_n}$  the radius of the regions around the  $n$  robots in the game (teammates and opponents). The case retrieval process consists in obtaining from the base the most similar case, the retrieved case. Therefore, it is necessary to compute the similarity between the current problem and the cases in the base. The similarity function indicates how similar a problem and a case are. In most cases, the function is defined by the distance between the ball and the robots in the problem and in the case.

$$Sim(p, c) = dist(B^c, B^p) + \sum_{i=1}^n dist(R_i^c, R_i^p), \quad (9)$$

where  $B^c$  is the position of the ball in the case and  $B^p$  its position in the problem,  $R_i^c$  the position of the Robot  $i$  in the case and  $R_i^p$  its position in the problem, and  $dist(a, b)$  is the gaussian distance between object  $a$  and  $b$ . This distance is computed as follows:

$$dist(a, b) = e^{-((a_x - b_x)^2 + (a_y - b_y)^2) / 2\tau^2}, \quad (10)$$

where  $\tau$  is the radius of the scope around the object. In this work,  $\tau$  is the same for the ball and robots positions. The Gaussian distance is used because the larger the distance between two points, the lower the similarity between them. Finally,  $\tau$  is used as a threshold that defines a maximum distance allowed for two points to have some degree of similarity: if  $dist(a, b) > \tau$ ,  $Sim(a, b) = 0$ .

Before a case can be reused, it might be necessary to adapt it to the present situation. Adaptation of a case means that the retrieved solution is modified, by translation, rotation or the addition of steps to the sequence of actions in the solution before it can be used. In this work, we assume that rotation and translation costs are small when compared to the cost of the additional steps, because the first two are trivial computations, while the performance of additional steps by the robots are actions that must be executed (in the simulator or in the real world), taking more time. Therefore, we define the cost as the number of steps added to the adapted solution. In this work, the case that will be reused is the one that maximizes the similarity while minimizing the adaptation cost.

In recent years, CBR has been used by several researchers in the Robotic Soccer domain. By far, the Robocup 2D Simulation League is the domain where most work has been done. To mention a few, Lin, Liu and Chen [19] presented a hybrid architecture for soccer players where the deliberative layer corresponds to a CBR system, Ahmadi *et al.* [2] presented a two-layered CBR system for prediction for the coach and Berger and Lämmel [5] proposed the use of a CBR system to decide whether a pass should be performed.

**Table 2.** The CB-HAMMQ algorithm.

---

```

Initialize  $\hat{Q}_t(s, a, o)$  and  $H_t(s, a, o)$  arbitrarily.
Repeat (for each episode):
  Initialize  $s$ .
  Repeat (for each step):
    Compute similarity and cost.
    If there is a case that can be reused:
      Retrieve and Adapt if necessary.
      Compute  $H_t(s, a, o)$  using Equation 5 with the
        actions suggested by the case selected.
    Select an action  $a$  using Equation 4.
    Execute the action  $a$ , observe  $r(s, a, o)$ ,  $s'$ .
    Update the values of  $Q(s, a, o)$  according to Equation 1.
     $s \leftarrow s'$ .
  Until  $s$  is terminal.
Until some stopping criterion is reached.

```

---

CBR has been also used in other Robocup Leagues. In the Small Size League, Srinivasan *et al.* [28] proposed a CBR planning for both offense and defense team behavior, for a team of two soccer playing robots; in the work by Marling *et al.* [23], CBR is used to help planning individual moves and team strategies. In the Four-Legged League, Karol *et al.* [17] presented high level planning strategies including a CBR system. Finally, the works of Ros *et al.* [26] presented the most ample use of CBR techniques in the Robotic Soccer domain, proposing the use of CBR techniques to handle retrieval, reuse and acquisition of a case base for the action selection problem of a team for the Four-Legged League. A more extensive review of the use of CBR in Robotic Soccer can be found in works by Burkhard and Berger [9] and by Ros [25].

## 4 Combining Case-Based Reasoning and Multiagent Reinforcement Learning

Bianchi, Ribeiro and Costa [6] state that there should be many methods that can be used to define a heuristic function for a HARL algorithm. For example, the same work makes use of information from the learning process itself to infer a heuristic at execution time, proposing a technique that derives a crude estimate of the transition probabilities, and then it propagates – from a final state – the correct policies which lead to that state. Bianchi, Ribeiro and Costa [7] employed prior domain knowledge to establish a very simple ad-hoc heuristic for speeding up learning in a Multiagent Reinforcement Learning domain.

In order to provide HAMRL algorithms with the capability of reusing previous knowledge from a domain, we propose a new algorithm, the Case-Based HAMMQ, that extends the HAMMQ algorithm, being capable of retrieving a case stored in a base, adapting it to the current situation, and building a heuristic function that corresponds to the case.

As the problem description  $P$  corresponds to one defined state of the set of states  $\mathcal{S}$  in an MDP, an algorithm that uses the RL loop can be implemented. Inside this loop, before action selection, we added steps to compute the similarity of the cases in the base with the current state and the cost of adaptation of these cases. A case is retrieved if the similarity is above a certain threshold, and the adaptation cost is low. After a case is retrieved, a heuristic is computed using Equation 5 and the actions suggested by the case selected. The complete CB-HAMMQ algorithm is presented in Table 2.

Although this is the first work that combines CBR with RL using

an explicit heuristic function, this is not the first work on combining the both fields. Drummond [11] was probably the first to use CBR to speed up RL, proposing to accelerate RL by transferring parts of previously learned solutions to a new problem. Sharma *et al.* [27] made use of CBR as a function approximator for RL, and RL as a revision algorithm for CBR in a hybrid architecture system; Juell and Paulson [15] exploited the use of RL to learn similarity metrics in response to feedback from the environment; Auslander *et al.* [3] used CBR to adapt quickly an RL agent to changing conditions of the environment by the use of previously stored policies and Li, Zonghai and Feng [18] proposed an algorithm that makes use of knowledge acquired by Reinforcement Learning to construct and extend a case base. Gabel and Riedmiller [12] makes uses of CBR to represent the learning function (the state value function  $V$ ) in RL, having an attribute-value based state/case representation and using K-Nearest Neighbor to predict the cases' solution. Using the same idea, these authors [13] extend an algorithm for multi-agent learning into a CBR framework, in an approach that makes easier the distributed learning of policies in cooperative multi-agent domains.

Our approach differs from all previous works combining CBR and MRL because of the heuristic use of the retrieved case. Bianchi, Ribeiro and Costa [7] proved that if the heuristic used is an admissible one, there will be a speed up in convergence time, if not, the use of the heuristic will not impede the RL method to converge to the optimal policy. As we use the case base as a heuristic, if the case base corresponds to an admissible heuristic there will be a speed up in the convergence time. But if the case base does not contain any useful case – or even if it contains cases that implement wrong solutions to the problem, the agent will learn the optimal solution anyway, by using the RL component of the algorithm [7]. Another difference of this proposal to previous works, such as the one presented in [8], is that a Multiagent RL algorithm is used, while others combined CBR with single-agent RL.

## 5 Experiments in the Robotic Soccer Domain

A set of empirical evaluations of the CB-HAMMQ approach were carried out in a proposed simulator for the robot soccer domain that extends the one proposed by Littman [20]. In this domain, called “Expanded Littman’s Soccer”, two teams, A and B, of three players each compete in a 10 x 15 grid presented in figure 1. Each team is composed by the goalie ( $g$ ), the defender ( $d$ ) and the attacker ( $a$ ). Each cell can be occupied by only one player. The actions that are allowed are: keep the agent still, move – north, south, east and west – or pass the ball to another agent. The action “pass the ball” from agent  $a_i$  to  $a_j$  is successful if there is no opponent in between them. If there is an opponent, it will catch the ball and the action will fail. Actions are taken in turns: all actions from one team’s agents are executed at the same instant, and then the opponents’ actions are executed. The ball is always with one of the players. When a player executes an action that would finish in a cell occupied by the opponent, it loses the ball and stays in the same cell. If an action taken by one agent leads it out the board, the agent stands still. When a player with the ball gets into the opponent’s goal, the trial ends and its team scores one point. The starting positions of all players are random, and the ball is given to one of the agents in a random fashion at the beginning of a trial.

To solve this problem, three algorithms were used: the Minimax-Q, described in section 2, the HAMMQ, described in section 2 and the CB-HAMMQ, proposed in section 4. Although this domain is still a simplified one, it is more complex than the original one pro-

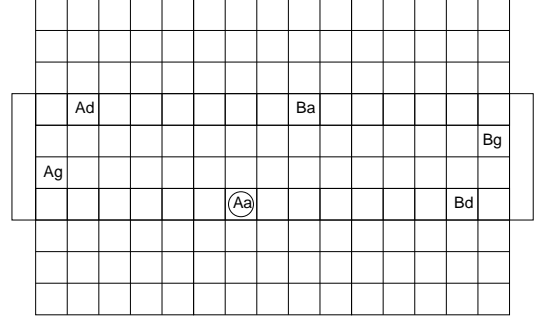


Figure 1. The “Expanded Littman’s Soccer” environment proposed.

posed by Littman: due to the size of the state space, it is not possible to use a lookup table containing all the states of the problem. In this work a variable resolution table similar to the one proposed by Munos and Moore [24] is used.

The heuristic used in the HAMMQ algorithm was defined using a simple rule: if holding the ball, go to the opponents’ goal, not taking into account the teammates’ and opponents’ positions, leaving tasks such as learning to pass the ball or to divert the opponent to the learning process.

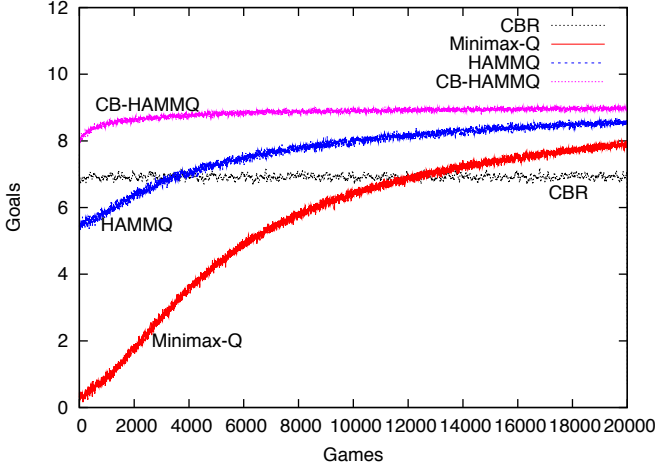
The heuristic value used in the CB-HAMMQ is computed during the games, as described in section 4. The case base used contains a set of basic cases that can be used without adaptation costs. The case base used in this experiment is composed of 5 basic cases, which cover the most significant situations that are observed during a game in the expanded Littman’s Soccer environment. These cases can be described as:

1. If the agent is with the ball and there is no opponent blocking it, then move to the goal.
2. If the agent is with the ball and there is an opponent blocking it, then move up.
3. If the agent is with the ball and there is an opponent blocking it, then move down.
4. If the agent is with the ball and a teammate is closer to the goal, then pass the ball to the other agent.
5. If the ball is with an opponent and the agent is close to the opponent, then stay in front of the opponent.

Is important to notice that this case base does not correspond to the optimal solution of the problem.

The reward the agents receive are the same for all algorithms: the agent that is holding the ball receives +100 every time it reaches the goal. This is a very simple reward scheme, but we decided to use it in this work to avoid the creation of a mismatch between the reward function used in training and the performance measure examined, which is the number of goals scored. Other reward schemes could be used, for example, one that gives rewards to intercepting the ball, losing the ball or correctly passing the ball, such as the one used by Kalyanakrishnan, Liu and Stone [16].

Thirty training sessions were run for the three algorithms, with each session consisting of 20,000 games of 10 trials. Figure 2 shows the learning curves for all algorithms when the learning team plays against an opponent moving randomly, and presents the average goal balance, which is the difference between goals scored and goals received by the learning team in each match. It is possible to verify that at the beginning of the learning phase Minimax-Q has worse performance than HAMMQ, and that this has a worse performance



**Figure 2.** Goals balance for the CBR, Minimax-Q, the HAMMQ and the CB-HAMMQ algorithms against a random opponent for the Expanded Littman's Robotic Soccer.

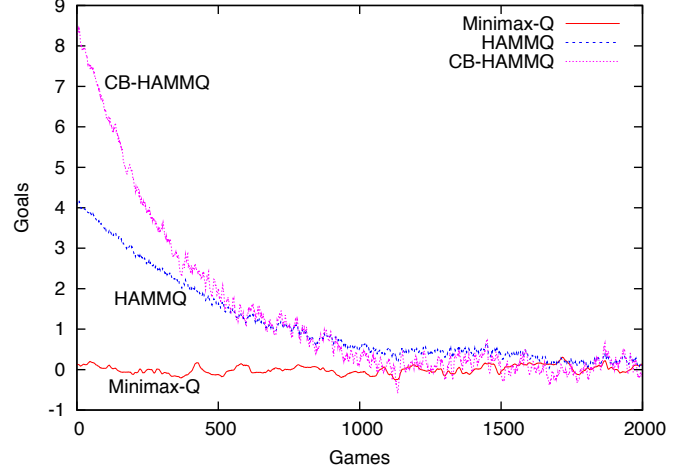
**Table 3.** Results for games against Random opponent.

| Algorithm                     | Goals made $\times$ goals conceded        |
|-------------------------------|---|
| Minimax-Q                     | $(140207 \pm 174) \times (38498 \pm 164)$ |
| HAMMQ                         | $(166208 \pm 150) \times (22065 \pm 153)$ |
| CB-HAMMQ                      | $(188168 \pm 155) \times (11292 \pm 140)$ |
| Games won $\times$ games lost |   |
| Minimax-Q                     | $(18297 \pm 33) \times (1037 \pm 28)$     |
| HAMMQ                         | $(19469 \pm 9) \times (27 \pm 4)$         |
| CB-HAMMQ                      | $(19997 \pm 1) \times (0 \pm 0)$          |

than CB-HAMMQ. As the matches proceed, the performance of the three algorithms become similar, as expected. As it can be seen in this figure, the Minimax-Q is still learning after 20,000 games: as it is slower than the other two algorithms, it will only reach the optimal solution after 100,000 games. In this figure the performance of a team of agents using only the case base can also be observed: a line with values close to 7. As the case base does not contain the optimal solution to the problem, the agents have a performance that is worse than the one presented by the other teams at the end of the learning process.

Figure 3 presents the learning curves ( the average goal balance at the end of a game) for the three algorithms when learning while playing against a learning opponent using Minimax-Q. It can be seen that CB-HAMMQ is better than HAMMQ and Minimax-Q at the beginning of the learning process. Student's  $t$ -test was used to verify the hypothesis that the use of heuristics speeds up the learning process. The result is that the CB-HAMMQ is better than HAMMQ until the 7,000<sup>th</sup> game when playing against a random opponent, and until the 500<sup>th</sup> game when playing against the Minimax-Q, with a level of confidence greater than 5%. The same test can be made comparing the CB-HAMMQ and the Minimax-Q: in this case, the first outperform the latter until the 20,000<sup>th</sup> game, when both are playing against a random opponent, and until the 1,000<sup>th</sup> game when the CB-HAMMQ is playing against the Minimax-Q. After these number of games the results of the algorithms are comparable, since the three algorithms converge to equilibrium.

Finally, Table 3 shows the average number of goals and the average number of games won at the end of 20,000 games while playing against a random opponent, and Table 4 presents the same data for games played against a Minimax-Q opponent, at the end of



**Figure 3.** Goals balance for Minimax-Q, the HAMMQ and the CB-HAMMQ algorithms against an opponent using Minimax-Q for the Expanded Littman's Robotic Soccer.

**Table 4.** Results for games against Minimax-Q opponent.

| Algorithm                     | Goals made $\times$ goals conceded      |
|-------------------------------|---|
| Minimax-Q                     | $(10299 \pm 234) \times (9933 \pm 240)$ |
| HAMMQ                         | $(10467 \pm 197) \times (9347 \pm 197)$ |
| CB-HAMMQ                      | $(11109 \pm 152) \times (8845 \pm 153)$ |
| Games won $\times$ games lost |   |
| Minimax-Q                     | $(848 \pm 60) \times (696 \pm 55)$      |
| HAMMQ                         | $(998 \pm 50) \times (530 \pm 43)$      |
| CB-HAMMQ                      | $(1145 \pm 37) \times (426 \pm 32)$     |

2,000 games. It can be seen in Table 4 that when Minimax-Q agents are playing against other Minimax-Q agents, the number of goals made and games won are approximately the same, while when CB-HAMMQ agents played against Minimax-Q ones, CB-HAMMQ team made more goals and won more games. CB-HAMMQ also won more games (1145, losing 425) and made more goals (11109) than the HAMMQ algorithm.

The parameters used in the experiments were the same for all the algorithms. The learning rate is  $\alpha = 0.9$ , the exploration/ exploitation rate was defined as being equal to 0.2 and the discount factor  $\gamma = 0.9$  (these parameters are similar to those used by Littman [20]). The value of  $\eta$  was set to 1. Values in the Q table were randomly initialized, with  $0 \leq Q(s_t, a_t, o_t) \leq 1$ .

## 6 Conclusion

This work presented a new algorithm, called Case-Based Heuristically Accelerated Minimax-Q (CB-HAMMQ), which allows the use of a case base to define heuristics to speed up the well-known Multiagent Reinforcement Learning algorithm Minimax-Q. This approach builds upon an emerging technique, the Heuristic Accelerated Reinforcement Multiagent Learning, in which MRL methods are accelerated by making use of heuristic information.

The experimental results obtained using a new domain proposed for the Robotic Soccer games showed that CB-HAMMQ attained better results than HAMMQ and Minimax-Q alone. For example, after playing 1000 learning trials against a random opponent (Figure 2), the Minimax-Q, still could not produce policies that scored many goals on the opponent, while the HAMMQ was able to score some goals but less than the CBR alone and the CB-HAMMQ. Another

interesting finding is that the number of goals scored by the CB-HAMMQ after 1000 trials was even higher than the number of goals scored by the CBR approach alone, indicating that the combination of the Reinforcement Learning and the case base out-performs the use of the case base on its own.

Finally, heuristic functions allow RL algorithms to solve problems where the convergence time is critical, as in many real time applications. Future works includes incorporating CBR in other well known Multiagent RL algorithms, like Minimax-SARSA [4], Minimax-Q( $\lambda$ ) [20] and expanding this framework to deal with General Sum Markov Games [20] using algorithms such as Nash-Q [14] and Friend-or-Foe Q-Learning [21]. Performing a game-theoretic analysis to determine if CB-HAMMQ is dominant against other strategies or if a mixed-strategy equilibrium is reached, using an approach based on [29], is also left as a future task.

## ACKNOWLEDGEMENTS

This work has been partially funded by the 2009-SGR-1434 grant of the Generalitat de Catalunya, the NEXT-CBR project, and FEDER funds. Reinaldo Bianchi acknowledge the support of the CNPq (Grants No. 201591/2007-3 and 453042/2010-4).

## References

- [1] Agnar Aamodt and Enric Plaza, 'Case-based reasoning: foundational issues, methodological variations, and system approaches', *AI Commun.*, 7(1), 39–59, (1994).
- [2] Mazda Ahmadi, Abolfazl Keighobadi Lamjiri, Mayssam M. Nevisi, Jafar Habibi, and Kambiz Badie, 'Using a two-layered case-based reasoning for prediction in soccer coach', in *Proc. of the Intern. Conf. of Machine Learning: Models, Technologies and Applications*, pp. 181–185. CSREA Press, (2003).
- [3] Bryan Auslander, Stephen Lee-Urban, Chad Hogg, and Héctor Muñoz-Avila, 'Recognizing the enemy: Combining reinforcement learning with strategy selection using case-based reasoning', in *Proceedings of the 9th European Conference on Case-Based Reasoning (ECCBR'08)*, pp. 59–73. Springer, (2008).
- [4] Bikramjit Banerjee, Sandip Sen, and Jing Peng, 'Fast concurrent reinforcement learners', in *Proceedings of the 17th International Joint Conference on Artificial Intelligence (IJCAI'01)*, Seattle, WA., pp. 825–832, (2001).
- [5] Ralf Berger and Gregor Lämmel, 'Exploiting past experience – case-based decision support for soccer agents', in *Proceedings of the 30th Annual German Conference on AI*, pp. 440–443. Springer, (2007).
- [6] Reinaldo A. C. Bianchi, Carlos H. C. Ribeiro, and Anna H. R. Costa, 'Accelerating autonomous learning by using heuristic selection of actions', *Journal of Heuristics*, 14(2), 135–168, (2008).
- [7] Reinaldo A. C. Bianchi, Carlos H. C. Ribeiro, and Anna Helena Reali Costa, 'Heuristic selection of actions in multiagent reinforcement learning', in *Proceedings of the 20th International Joint Conference on Artificial Intelligence (IJCAI 2007)*, Hyderabad, India, pp. 690–695, (2007).
- [8] Reinaldo A. C. Bianchi, Raquel Ros, and Ramon Lopez De Mantaras, 'Improving reinforcement learning by using case based heuristics', in *Proceedings of the 8th International Conference on Case-Based Reasoning (ICCBR'09)*, eds., Lorraine McGinty and David C. Wilson, volume 5650 of *Lecture Notes in Artificial Intelligence*, pp. 75–89, Springer, (2009).
- [9] Hans-Dieter Burkhard and Ralf Berger, 'Cases in robotic soccer', in *Proceedings of the 7th International Conference on Case-Based Reasoning (ECCBR'07)*, pp. 1–15. Springer, (2007).
- [10] Luiz A. Celiberto, Carlos H. C. Ribeiro, Anna Helena Reali Costa, and Reinaldo A. C. Bianchi, 'Heuristic reinforcement learning applied to robocup simulation agents', in *RoboCup*, eds., Ubbo Visser, Fernando Ribeiro, Takeshi Ohashi, and Frank Dellaert, volume 5001 of *Lecture Notes in Computer Science*, pp. 220–227. Springer, (2007).
- [11] Chris Drummond, 'Accelerating reinforcement learning by composing solutions of automatically identified subtasks', *Journal of Artificial Intelligence Research*, 16, 59–104, (2002).
- [12] Thomas Gabel and Martin Riedmiller, 'CBR for state value function approximation in reinforcement learning', in *In Proceedings of the 6th International Conference on Case Based Reasoning (ICCBR 2005)*, pp. 206–221. Springer, (2005).
- [13] Thomas Gabel and Martin Riedmiller, 'Multi-agent case-based reasoning for cooperative reinforcement learners', in *Proceedings of the 8th European Conference on Case-Based Reasoning (ECCBR '06)*, pp. 32–46. Springer, (2006).
- [14] Junling Hu and Michael P. Wellman, 'Nash Q-learning for general-sum stochastic games', *Journal of Machine Learning Research*, 4, 1039–1069, (2003).
- [15] Paul Juell and Patrick Paulson, 'Using reinforcement learning for similarity assessment in case-based systems', *IEEE Intelligent Systems*, 18(4), 60–67, (2003).
- [16] Shivaram Kalyanakrishnan, Yaxin Liu, and Peter Stone, 'Half field offense in robocup soccer: A multiagent reinforcement learning case study', in *RoboCup*, eds., Ubbo Visser, Fernando Ribeiro, Takeshi Ohashi, and Frank Dellaert, volume 5001 of *Lecture Notes in Computer Science*, pp. 72–85. Springer, (2007).
- [17] Alankar Karol, Bernhard Nebel, Christopher Stanton, and Mary-Anne Williams, 'Case based game play in the robocup four-legged league part i: the theoretical model', in *RoboCup*, eds., Daniel Polani, Brett Browning, Andrea Bonarini, and Kazuo Yoshida, volume 3020 of *Lecture Notes in Computer Science*, pp. 739–747. Springer, (2003).
- [18] Yang Li, Chen Zonghai, and Chen Feng, 'A case-based reinforcement learning for probe robot path planning', in *Proceedings of the 4th World Congress on Intelligent Control and Automation, Shanghai, China*, pp. 1161–1165, (2002).
- [19] Yi-Sheng Lin, Alan Liu, Kuan-Yu Chen, 'A hybrid architecture of case-based reasoning and fuzzy behavioral control applied to robot soccer', in *Workshop on Artificial Intelligence, International Computer Symposium (ICS2002)*, Hualien, Taiwan, National Dong Hwa University, (2002).
- [20] Michael L. Littman, 'Markov games as a framework for multi-agent reinforcement learning', in *Proceedings of the 11th International Conference on Machine Learning (ICML'94)*, pp. 157–163, (1994).
- [21] Michael L. Littman, 'Friend-or-foe q-learning in general-sum games', in *Proceedings of the 18th International Conference on Machine Learning (ICML'01)*, pp. 322–328. Morgan Kaufmann, (2001).
- [22] Ramon López de Mantaras, David McSherry, Derek Bridge, David Leake, Barry Smyth, Susan Craw, Boi Faltings, Mary Lou Maher, Michael T. Cox, Kenneth Forbus, Mark Keane, Agnar Aamodt, and Ian Watson, 'Retrieval, reuse, revision and retention in case-based reasoning', *Knowl. Eng. Rev.*, 20(3), 215–240, (2005).
- [23] Cynthia Marling, Mark Tomko, Matthew Gillen, David Alexander, and David Chelberg, 'Case-based reasoning for planning and world modeling in the robocup small size league', in *IJCAI-03 Workshop on Issues in Designing Physical Agents for Dynamic Real-Time Environments*, (2003).
- [24] Remi Munos and Andrew Moore, 'Variable resolution discretization in optimal control', *Machine Learning*, 49(2/3), 291–323, (2002).
- [25] Raquel Ros, *Action Selection in Cooperative Robot Soccer using Case-Based Reasoning*, Ph.D. dissertation, Universitat Autònoma de Barcelona, Barcelona, 2008.
- [26] Raquel Ros, Josep Lluís Arcos, Ramon López de Mantaras, and Manuela Veloso, 'A case-based approach for coordinated action selection in robot soccer', *Artificial Intelligence*, 173(9-10), 1014–1039, (2009).
- [27] Manu Sharma, Michael Holmes, Juan Carlos Santamaría, Arya Irani, Charles Lee Isbell Jr., and Ashwin Ram, 'Transfer learning in real-time strategy games using hybrid CBR/RL', in *Proceedings of the 20th International Joint Conference on Artificial Intelligence (IJCAI 2007)*, Hyderabad, India, pp. 1041–1046, (2007).
- [28] Thanukrishnan Srinivasan, K. Aarthi, S. Aishwarya Meenakshi, and M. Kausalya, 'CBRobosoc: An efficient planning strategy for robotic soccer using Case-Based Reasoning', in *Proceedings of the International Conference on Computational Intelligence for Modelling Control and Automation and International Conference on Intelligent Agents Web Technologies and International Commerce (CIMCA '06)*, p. 113. IEEE Computer Society, (2006).
- [29] William E. Walsh, Rajarshi Das, Gerald Tesauero, and Jeffrey O. Kephart, 'Analyzing complex strategic interactions in multi-agent games', in *AAAI-02 Workshop on Game Theoretic and Decision Theoretic Agents*, pp. 109–118, (2002).